

# Inferring the T-cells repertoire dynamics of healthy individuals

Meriem Bensouda Koraichi,<sup>1</sup> Silvia Ferri,<sup>1,2</sup> Aleksandra M Walczak,<sup>1,\*</sup> and Thierry Mora<sup>1,\*</sup>

<sup>1</sup>*Laboratoire de physique de l'École normale supérieure, CNRS, PSL University, Sorbonne Université, and Université Paris Cité, 75005 Paris, France*

<sup>2</sup>*Dipartimento di Fisica e Astronomia dell'Università di Bologna, 40126 Bologna, Italy*

The adaptive immune system is a diverse ecosystem that responds to pathogens by selecting cells with specific receptors. While clonal expansion in response to particular immune challenges has been extensively studied, we do not know the neutral dynamics that drive the immune system in absence of strong stimuli. Here we learn the parameters that underlie the clonal dynamics of the T-cell repertoire in healthy individuals of different ages, by applying Bayesian inference to longitudinal immune repertoire sequencing (RepSeq) data. Quantifying the experimental noise accurately for a given RepSeq technique allows us to disentangle real changes in clonal frequencies from noise. We find that the data are consistent with clone sizes following a geometric Brownian motion, and show that its predicted steady state is in quantitative agreement with the observed power-law behaviour of the clone-size distribution. The inferred turnover time scale of the repertoire increases substantially with patient age, and depends on the clone size in some individuals.

## I. INTRODUCTION

The adaptive immune system protects us from many infections including those caused by pathogenic challenges that did not exist when we were born. This amazing plasticity is encoded, in part, in a diverse repertoire of T cells carrying surface receptors capable of recognizing different antigens, which trigger an immune response. About  $10^8$  new T cells are estimated to be generated and enter the periphery in human adults every day [1, 2], where they undergo specific proliferation due to antigen stimulation but also non-specific divisions [3, 4] and death. These processes together result in clone sizes of different T cells that differ over a few orders of magnitude, forming long tailed distributions [5, 6]. The total number of different T cell clones is estimated between  $10^8$  and  $10^{10}$  [7–9]. Qualitatively describing the T cell clonal dynamics in the periphery is important for predicting long- as well as short-term immune response and to understand the maintenance of immune memory.

A lot of effort has been put into describing antigen specific response and memory formation [10–12]. At any given timepoint the majority of the T cell repertoire is not always directly involved in fighting the current antigenic challenge. Yet, processes such as homeostasis [3] and unspecific signals in both naive and memory subrepertoires result in frequency changes of background clones. Many first-principles models of naive T cells dynamics have been proposed to study the balance between thymic output and peripheral proliferation and death [2, 4, 13–15]. The role of competition for antigens between T cells has been pointed out [16], as well the effect of cross-reactivity [17] (the ability for one T cell to recognize different antigens), the relative size of a primary versus secondary response to similar antigens [18], or the effect of beneficial mutations [19]. These studies highlight

the importance for the naive repertoire of clonal expansions that are not necessarily linked to specific challenges. While these models were instrumental in advancing our understanding of bulk repertoire dynamics, and allowed for the interpretation of deuterated water and bromium staining experiments that describe cellular lifetimes [20], the class of models that are consistent with the data is still large and unexplored.

Thanks to advances in immune repertoire sequencing (RepSeq) [21–23], dynamical models can now be assessed directly against repertoire data at the clonal level. RepSeq experiments isolate and sequence the T cell receptors (TCR) in a blood sample of individuals. By counting reads with the same TCR sequence, one can estimate the frequency of the corresponding clone (defined as the set of cells with the same receptor) in blood. Even single repertoire snapshots can be informative about the dynamics: the distribution of clone sizes follows a power law [6, 24–27], in accord with proposed models of stochastic growth and death [5]. Taking samples from the same individual at different timepoints allows for tracking the evolution of TCR clone sizes in time. The longitudinal experiments that have been performed in healthy donors [28, 29] suggest that the repertoire is relatively stable over years.

Our main goal in this article is to characterize the dynamics of the unstimulated background repertoire. We use an inverse approach to learn models of stochastic TCR clonal dynamics directly from data, collecting human TCR RepSeq datasets where we could identify at least two time points between which there was no reported specific acute antigenic stimulation [28–32]. A key aspect of our method is the treatment of experimental noise, which confounds naive analyses of stochastic time traces. The method first quantifies both the sampling and natural biological noise thanks to replicate RepSeq experiments [33, 34], and then infers the parameters of a stochastic dynamical model to describe the trajectories of each TCR clone population in a healthy individual. We explicitly show how correcting for noise allows us to

---

\*Corresponding authors. These authors contributed equally.

robustly learn the underlying dynamics.

A recent study [35] has investigated the formation of the T cell receptor during development and its maintenance into adulthood. Here we focus on healthy adult repertoires that are already shaped during the first years of an individual's life, and ask how they evolve and get renewed. We extract clone turnover time scales, and describe how these time scales depend on the person's age. Characterizing this baseline dynamics is an important step towards interpreting TCR dynamics in the presence of antigenic stimuli.

## II. RESULTS

### A. Longitudinal sampling of TCR repertoires of healthy individuals

T-cell repertoires are large ecosystems in which each species is a clone of T cells carrying the same TCR  $i$  formed by a unique pair of  $\alpha$  and  $\beta$  chains. The dynamics of this system is characterized by the time course of the number of cells carrying each receptor,  $n_i(t)$ . This number can be accessed indirectly through TCR repertoire sequencing (RepSeq), obtained by sequencing the TCR of small samples of peripheral blood mononuclear cells (PBMC), giving us a read count  $\hat{n}_i(t)$  for a given chain at different timepoints (Fig. 1A). Because the two chains are not paired in the data, from here on we define clones as collections of cells having the same  $\alpha$  or  $\beta$  chain, which we will refer to as clonotypes. This approximation is justified by the low occurrence of TCRs that share one chain but not the other [32].

We collected repertoire data from 9 individuals P1-P9, aged 18-57, sampled at various time points from one month up to 3 years apart, with and without biological replicates.  $\beta$  chain repertoires were sequenced for all samples, and  $\alpha$  chains only for individual P6. The properties of the datasets, including their number of clones  $N_c$ , total read counts  $N_r = \sum_{i=1}^{N_c} \hat{n}_i$ , age, library preparation (from genomic DNA or from mRNA), and chain, are summarized in Fig. 1B and in Table S1. Because P3, P4, P6, and P9 were included in a vaccination study, they had received a shot of the YFV 17D yellow fever vaccine (P3, P4, P6) or of the influenza vaccine (P9) 45 days prior to the first time point, after the decay of their T-cell response, so we assume that the dynamics of vaccine-specific T does not affect much our analysis of the global repertoire.

A major challenge when analyzing RepSeq data is that the measured abundances  $\hat{n}_i(t)$  only provide a noisy reflection of the true ones  $n_i(t)$ . Observed differences between datasets thus result from a combination of the repertoire dynamics and biological and experimental noise. The magnitude of that noise can be assessed by comparing the normalized clonotype frequencies  $\hat{f} = \frac{\hat{n}_i}{N_r}$  between two biological replicates obtained at the same time point in the same individual (Fig. 1C, blue dots). By

contrast, comparing those frequencies between two timepoints separated by one year (Fig. 1C, red dots) show a larger dispersion, and a slight overall decrease of clonotype frequencies. Our goal is to measure this difference quantitatively.

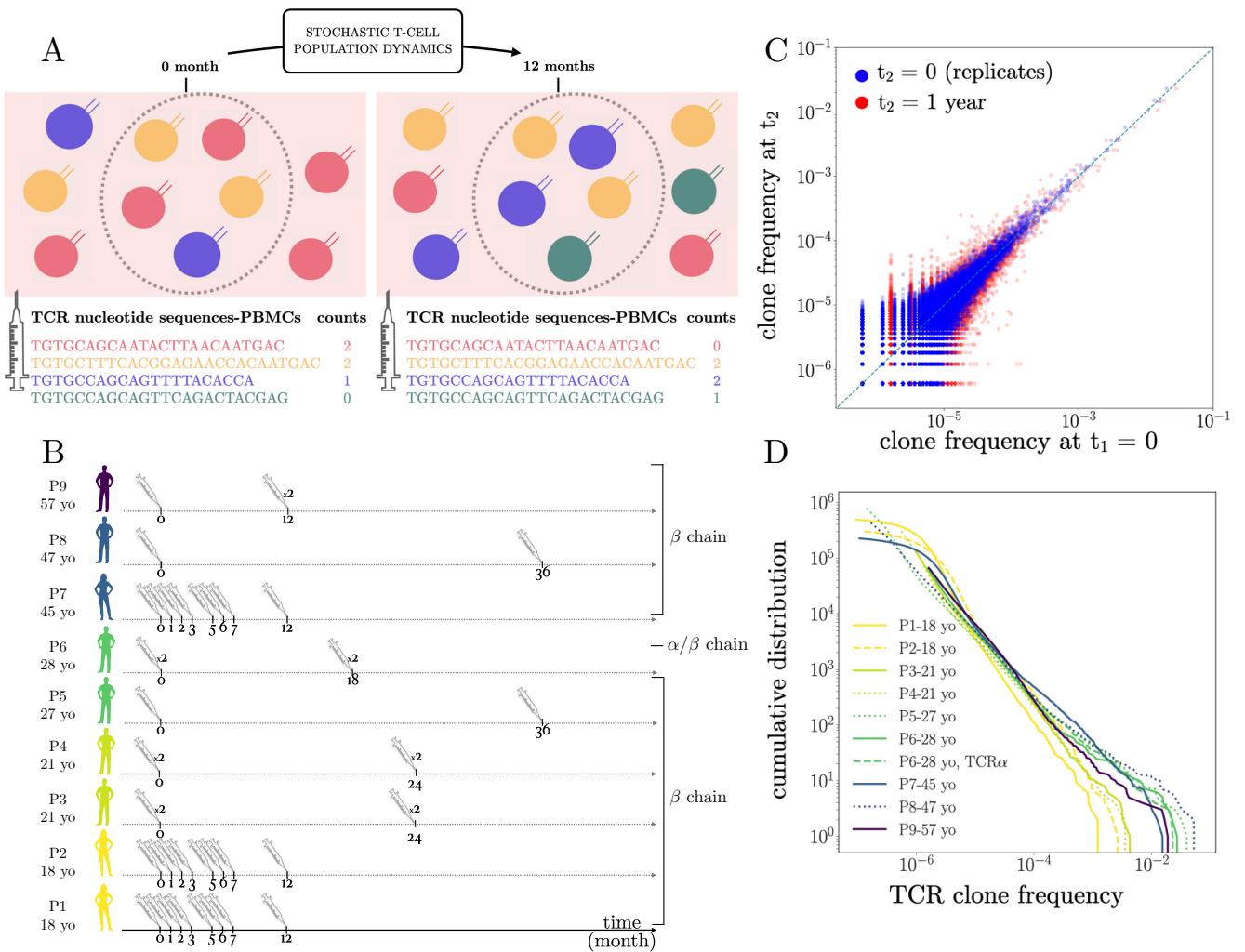
Another difficulty arises from the observation that clonotype frequencies are highly heterogeneous, with their distribution following a power law  $P(\hat{n}) \propto \hat{n}^{-1-\alpha}$  spanning no less than 4 orders of magnitude, with an exponent  $\alpha \approx 1$  which is largely invariant across individuals and timepoints (Fig. 1D), as previously reported [5, 6, 35]. This implies that most clonotypes have very low abundance and are thus particularly subject to sampling and experimental noise.

### B. Mathematical model of stochastic clonal dynamics

The dynamics of T cell clones is driven by the proliferation and death of cells belonging to them. In addition, new clones with their distinct TCR are continually produced and released by the thymus, although the rate of thymic exports decays rapidly with age [1]. Cell division, death, and introduction of new clones constitute the basis of our model (Fig. 2A). Cell division may be caused by antigen stimulation (both self and foreign) or by cytokine and growth factors, and cells die by lack of stimulation or by apoptotic signals. Even in the absence of strong and chronic antigenic stimuli, T cell clonotype abundances display stochastic trajectories due to either weak stimulation, repertoire homeostasis and demographic fluctuations. In addition, individuals may get mild infections over the course of months and years. Since these events are numerous and unknown, we model them by an effectively random net growth rate (divisions minus deaths). It can be shown [5, 36] that on time scales much longer than the typical resolution time of infections, each clonotype size  $n_i(t)$  may then be modeled by a geometric Brownian motion (GBM). Its evolution is governed by an effective mean net growth, to which random fluctuations are added to account for bursts of proliferation and decay:

$$\frac{d\ln n_i(t)}{dt} = -\tau^{-1} + \theta^{-1/2} \eta_i(t), \quad (1)$$

where  $\eta_i(t)$  is a clonotype-specific white noise of zero mean and unit amplitude. Note that the mean growth rate of clones,  $-\tau^{-1}$ , is typically negative. On average, each clone should decay to make room for new thymic exports, because of homeostatic pressures that control the total number of cells. In this interpretation,  $\tau$  is the typical decay and turnover time of each clone, which would evolve with time as  $n_i(t) = n_i(0)e^{-t/\tau}$  in the absence of fluctuations. But recall that this is just an average—many clones do not decay, but instead undergo episodes of large growth and decay, as illustrated by simulations



**FIG. 1: Longitudinal tracking of T-cell repertoires.** **A.** Experimental workflow. PBMC from a healthy individual are extracted at two timepoints, and their TCR repertoire sequenced, yielding lists of clonotypes with count numbers corresponding to the number of individual measurements (or reads). The way in which the two sampled repertoires has changed between the two timepoints is predicted by a stochastic model of the dynamics of T cell clones. **B.** Summary of the TCR  $\alpha$  and  $\beta$  repertoire data used in this study. 9 individuals from 5 studies, aged 18-57, male and female, were included. When available, replicate experiments are annotated with  $\times 2$ . Datasets were produced using two different sequencing technologies based on cDNA and gDNA. **C.** Typical scatter plot of frequencies of TCR clones in two samples from the same individual P9. Blue: two biological replicates obtained on the same day show the effect of experimental noise. Red: 2 samples taken 1 year apart show a larger spread, resulting from a combination of real changes and noise. The goal of the analysis is to disentangle real changes from the noise. **D.** Cumulative distributions of TCR frequencies, which follow a universal power law in all samples and donors, with exponent  $\approx 1$ .

of (1) in Fig. 2B. The typical amplitude of these fluctuations grows with time as  $\sqrt{t/\theta}$  (dashed lines). Thus,  $\theta$  may be interpreted as the typical time it takes for a clone to rise or decay above or below the typical behaviour by one log-unit.

In addition to being biological motivated, the proposed dynamics have the desirable property that, in the presence of a constant rate of thymic exports, the distribution of clone sizes is predicted to evolve in time towards a perfect power law,  $P(n) \propto n^{-1-\alpha}$ , with exponent  $\alpha = 2\theta/\tau$  given by twice the ratio of the two time scales of the

model [5]. This is illustrated in Fig. 2C on simulated repertoires at steady state, and agrees well with the empirical distributions of Fig. 1D.

Our goal is to capture the parameters of these dynamics that is informative about the repertoire turnover timescales, while constraining the experimentally observed clone size frequency distribution. Our approach assumes that on the timescales of the analysis we do not observe signals of strong and specific antigenic stimulation. It also ignores potential dependences on the size of the clone, which could be mediated by phenotypic dif-

ferences between clones. This last assumption will be revisited later.

### C. Model inference

We estimate the parameters  $(\theta, \tau)$  of the dynamics in Eq. (1) from the observed clonotype abundance trajectories using a Bayesian approach for the posterior distribution of parameters given the data:

$$(\tau^*, \theta^*) = \arg \max_{\tau, \theta} \prod_{i=1}^{N_c} P(\hat{n}_i(t_1), \hat{n}_i(t_2) | \tau, \theta), \quad (2)$$

where  $t_1$  and  $t_2$  are the times of the two samples.

We use two methods to learn the model parameters: *naive inference* and *full inference*. The naive inference assumes the empirical abundances faithfully represents the real clonal abundances  $n_i$  through a simple proportionality rule,  $\hat{n}_i \approx (N_r/N_{\text{cell}})n_i$ , where  $N_{\text{cell}} = \sum_i n_i$  is the total number of T cells in the body. In practice, we work with clonotype frequencies  $f_i = \hat{n}_i/N_r$ , and  $f_i = n_i/N_{\text{cell}}$ , so that this assumption becomes  $\hat{f}_i = f_i$ . Further assuming that the total number of cells  $N_{\text{cell}}$  is approximately constant in time at steady state,  $\hat{f}_i$  is then governed by the same equation (1) as  $n_i$ . We take advantage of the closed solution available for the propagator associated to the GBM, (4), to maximize the log-likelihood (see Methods). This maximization is equivalent to plotting a histogram of the change in log-frequencies between the two timepoints, and simply read off  $\tau^{-1}$  and  $\theta^{-1}$  as the negative mean and the variance of the distribution divided by  $t = t_2 - t_1$  (Fig. 3A), consistent with their biological interpretation.

The *full inference* incorporates the fact that the observed clonotype abundances are contaminated by biological (mRNA expression) and experimental noise sources (sequencing errors, stochastic PCR amplification, and sampling), which means they do not correspond exactly to the clonotype abundances. To give a sense of just the sampling issue, a PBMC sample of  $\sim 1$  mL contains about 1 million cells, yielding about 1 million reads. By comparison, the organism contains of the order of  $10^{11}$  T cells. TCR clonotype frequencies are thus extrapolated from observing a fraction  $10^6/10^{11} \approx 10^{-5}$ , or 0.001%, of the whole repertoire [9]. In addition, not all cells are captured, and each cell may be represented by multiple reads, either through sequencing of multiple mRNA from the same cell, or from PCR amplification, depending on the context. To address these sources of uncertainty, in the full inference approach we introduce an error model [33] relating observed frequencies  $\hat{f}$  to their true value  $f$  probabilistically through the transfer function  $P(\hat{f}|f)$  (Fig. 3B). We use the previously introduced a software tool, NoisET [34], which learns such a noise model from replicate RepSeq experiments (see Methods).

We applied NoisET to individuals P3, P4, P6 and P9 for whom replicates were available. The noise model as-

sumes that the read count  $\hat{n}$  of each clone is drawn from a negative binomial distribution, whose variance grows with the frequency as  $\text{Var}(\hat{n}) = fN_r + a(fN_r)^b$ , with two learnable parameters  $a, b$ . In addition, since true frequencies are unknown, we assume as a prior that frequencies are distributed according to a power law  $\rho(f) \propto f^{-1-\alpha}$  with a cut-off  $f > f_{\min}$ , with  $\alpha$  and  $f_{\min}$  another two parameters. These parameters are reported for all individuals and time points in Fig. S1.

Once the noise model has been learned using NoisET, the likelihood of the data is computed by summing over the latent variables  $f_1 = f_i(t_1)$  and  $f_2 = f_i(t_2)$ :

$$P(\hat{n}_i(t_1), \hat{n}_i(t_2) | \tau, \theta) = \iint_{f_{\min}}^1 df_1 df_2 \rho(f_1) P(f_2 | f_1, \tau, \theta) \\ \times P(\hat{f}_1 | f_1) P(\hat{f}_2 | f_2), \quad (3)$$

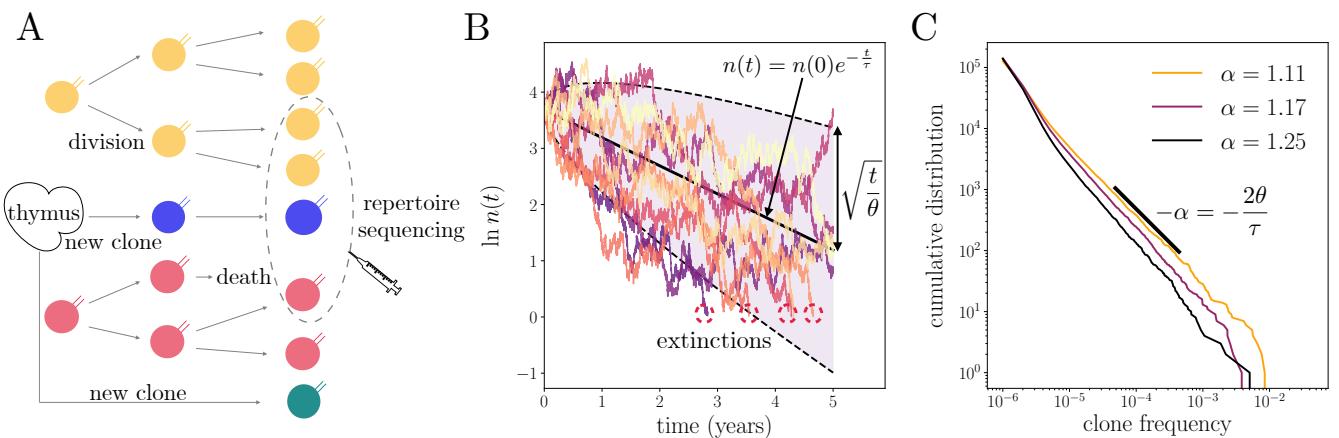
where  $P(f_2 | f_1; \tau, \theta)$  is the propagator of the geometric Brownian motion Eq. 1, and  $\hat{f}_j = \hat{n}_i(t_j)/N_r(t_j)$ ,  $j = 1, 2$ .

To explore the dependence of the  $\tau$  and  $\theta$  parameters on the frequency of clonotypes, and to eliminate clones that are not seen at both timepoints, we can generalize the formulas above to include only clonotypes with frequencies larger than a specific threshold  $\hat{f}_{\text{th}}$ , which modifies the normalization of the maximum likelihood estimator (see Methods).

### D. Validation of the inference methods on synthetic data

We first test the naive and full model inference on simulated RepSeq samples. We simulate  $10^{10}$  cells corresponding to  $\sim 10^8$  synthetic longitudinal trajectories designed to mimic as closely as possible the features of the real repertoire data at time points two years apart. The initial size  $n_i(t_1)$  of each clone is drawn from the steady state distribution of the GBM with a constant source (see Methods). Then Eq. 1 is simulated between times  $t_1$  and  $t_2$  with an extinction condition when  $n_i < 1$ , and with a source of new clones whose rate of introduction is matched to the mean extinction rate (see Methods). We varied the two timescales of the model,  $\tau$  and  $\theta$ , from months to years, while keeping  $\alpha = 2\theta/\tau$  within the observed experimental range 1.1–1.25 [35].

We model experimental sampling using a negative binomial distribution with variance parameters  $a = 0.7$  and  $b = 1.1$ . Sequencing depth was set to  $N_r = 10^6$  reads at both time points (we checked that asymmetric numbers of reads at each timepoint did not affect the results, see Fig. S2), resulting in  $\sim 10^5$  sampled distinct clonotypes. For each set of parameters we generated 10 longitudinal datasets to assess errors. We then performed the naive and full inference methods on these datasets, restricted to clones with  $\hat{f}_1 \geq \hat{f}_{\text{th}}$  and  $\hat{f}_2 > 0$ , and compared the inferred values of  $\tau$  and  $\theta$  to the true ones (Fig. 3C-D).



**FIG. 2: Stochastic model of repertoire dynamics.** **A.** T cells are introduced in the peripheral immune system by thymic export, providing a source of new TCR clones. T cells belonging to a specific TCR clone (labeled by their color) divide and die depending on their interactions with the antigenic environment, increasing or reducing the abundance of its TCR in the repertoire. This process is modeled by a geometric Brownian motion **B.** Example traces of TCR abundances simulated from the model Eq. 1 with  $n(0) = 40$ , with  $\tau = 2$  years and  $\theta = 1.11$  year. Clones that reach abundance  $< 1$  go extinct (red circles). The typical trend is for clones to decay exponentially with time scale  $\tau$  (black solid line). Stochastic events of growth and decay account for a broad variability of individual traces, whose magnitude grows as  $\sqrt{t/\theta}$  with time (shaded area) in logarithmic scale. **C.** Cumulative frequencies distributions of synthetic TCR clone abundances. The model predicts a power law of exponent  $\alpha = 2\theta/\tau$ . Different values of  $\tau$  and  $\theta$  were used to lead to different values of the exponent  $\alpha$ . Parameters:  $\tau = 2$  years,  $N_{\text{cell}} = 10^{10}$ ,  $n_0 = 40$ .

While the full inference (blue points) works for all values of the parameters and frequency threshold, the naive inference (red points) performs poorly for large values of the parameters. Increasing the cutoff frequency to  $f_{\text{th}} = 10^{-4}$  improves the naive inference by limiting the effect of the sampling noise, which is relatively smaller in large clones. For lower values of the threshold, the more numerous small clones dominate the inference, yielding an erroneous estimate. However, since the naive inference does not require replicates or a noise model, and is faster to implement, it provides a practical solution for learning  $\tau$  and  $\theta$  for large clones.

### E. Analysis of repertoires

We applied the full inference to longitudinal data sets of healthy individual TCR repertoires presented in Fig. 1 for which replicates were available, focusing on large enough clones ( $\hat{f}_1 \geq f_{\text{th}} = 10^{-5}$ ). With this cutoff we limit experimental noise and focus mainly on memory clones, since large clones are more likely to have arisen from expansion and belong to the memory pool [14]. For all individuals, the inferred values of  $\tau$  and  $\theta$  fall close to a line defined by  $\tau \approx 2\theta$  (Fig. 4A) corresponding to a predicted exponent of  $\alpha \approx 2\theta/\tau = 1$  in the power law of the clone size distribution. This result is in agreement with empirical observations of Fig. 1D. A more refined comparison of the predicted exponent,  $2\theta/\tau$ , with the one directly inferred from the distribution of clone sizes,

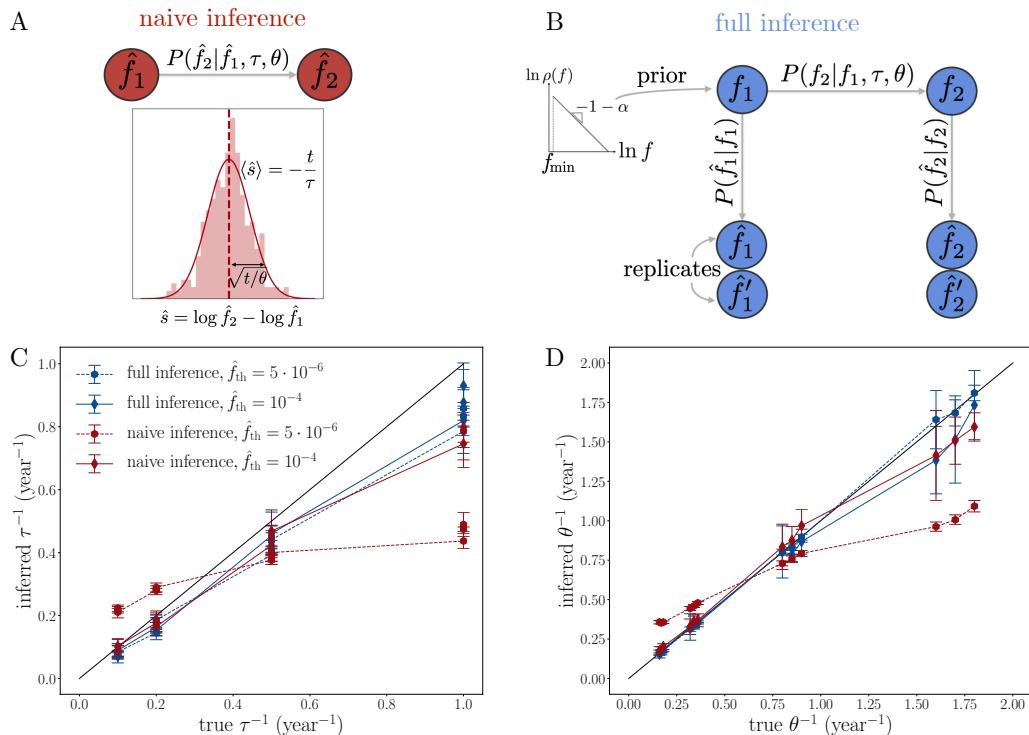
$\alpha$ , gives consistent but noisy results (Fig. 4B), primarily because of the narrow range of values of  $\alpha$  (0.9–1.2) and the small number of individuals. We note that the two inferred values of  $\alpha$  use completely independent pieces of information, namely the clone size distribution in one case, and the dynamics of clone sizes in the other.

Since our approach is probabilistic, it provides as a byproduct the posterior distribution of the fold change of individual clones (see Methods). The average of this posterior over clones agrees very well with the model propagator (4) (Fig. S3), validating its consistency with the data.

The turnover time  $\tau$  increases sharply with age, from a few years at age 21 to  $\sim 50$  years at age 57 (Fig. 4C). Since the ratio  $2\theta/\tau$  is constrained to be  $\approx 1$ , this implies that the amplitude of the stochastic stimulations,  $\theta^{-1}$ , decreases with age. The TCR repertoire is more dynamic, with faster turnover, for young individuals, who also have a larger rate of introduction of new TCR clones from the thymus than older individuals. At the same time, a turnover time of  $\sim 20$  years at the age of 40 suggests that the repertoires of adults remain dynamic despite greatly reduced thymic output.

For individual P6, both TCR $\alpha$  and TCR $\beta$  RepSeq samples were available. We recover very similar dynamic parameters for both receptor chains (Fig. 4A-B). This justifies our hypothesis that the bulk sequencing of single chains captures well the dynamics of  $\alpha/\beta$  clonotypes.

For comparison, we also applied the naive inference procedure, which allows us to include all 9 patients even



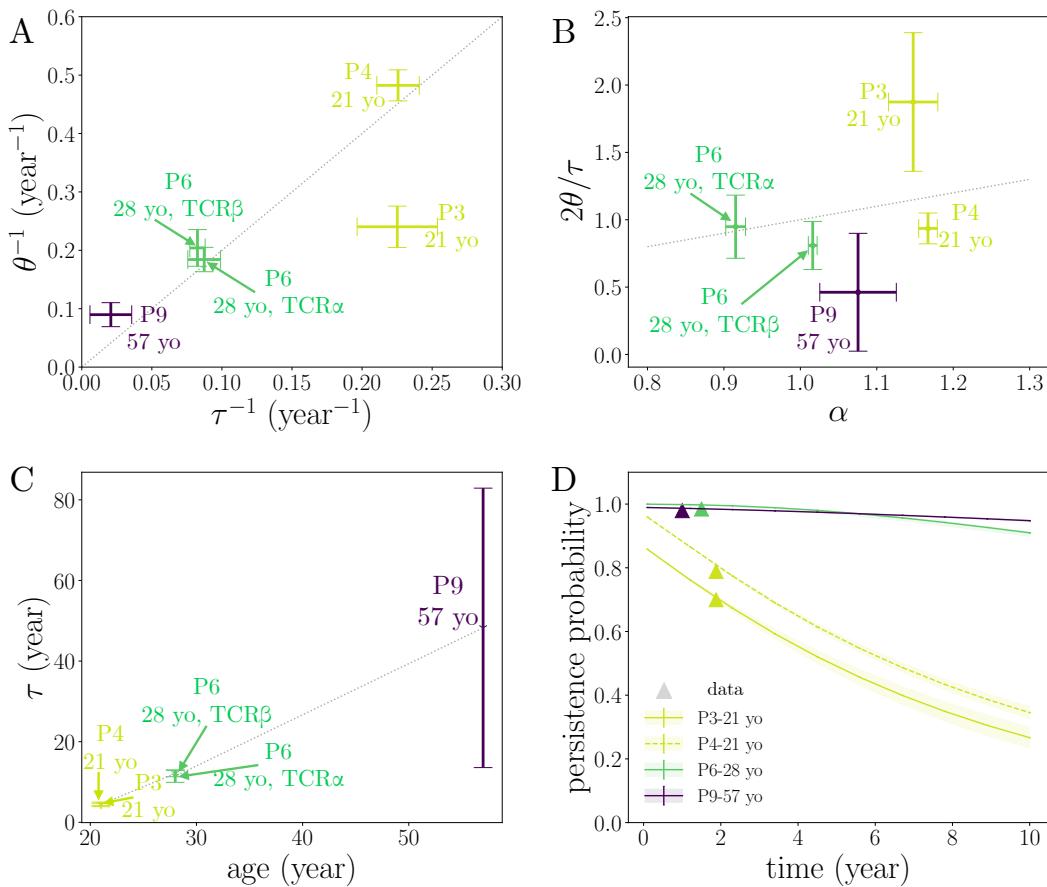
**FIG. 3: Inferring dynamical parameters from data.** **A.** Naive inference. The empirical clonotype frequencies in the RepSeq samples at the two time points,  $\hat{f}_1$  and  $\hat{f}_2$ , are treated as the true ones,  $f_1$  and  $f_2$ . We estimate the two parameters  $\tau$  and  $\theta$  from the average and standard deviation over all observed TCR clones of the log-fold change in frequency between two-time points,  $\log(\hat{f}_2/\hat{f}_1)$ , which the model predicts is distributed normally. The distribution of  $\hat{s}$  is shown in light red, and the Gaussian fit in solid red. **B.** Full inference. The empirical frequencies  $\hat{f}$  are modeled as noisy read-outs of the true ones  $f$ , through a probabilistic noise model. First, the noise model  $P(\hat{f}|f)$  is inferred from replicate experiments such as shown in Fig. 1C. The inference procedure also learns the distribution of frequencies  $\rho(f)$ , assumed to follow a power law with adjustable exponent  $\alpha$  and minimal frequency  $f_{\min}$ . Second, using the noise model, the parameters of the dynamical propagator  $P(f_2|f_1, \tau, \theta)$  are inferred from two timepoints, where  $f_1$  and  $f_2$  are treated as latent variables and  $\hat{f}_1$ ,  $\hat{f}_2$  as observables, using a Maximum likelihood estimator. **C-D.** Validation of naive and full inference models on synthetic data. Model parameters:  $t_2 - t_1 = 2$  years,  $\tau = 1, 2, 5, 10$  years,  $\alpha = 1.11, 1.18, 1.25$ , with all 12 combinations tested; number of cells  $N_c = 10^{10}$ ; initial clone size  $n_0 = 40$ ; the other parameters (number of clones, thymic output) are deduced assuming steady state (see Methods). Sampling model: number of sampled reads  $N_r = 10^6$ ; noise model parameters  $a = 0.7$  and  $b = 1.1$ . Error bars are standard deviations over 10 simulations.

when replicates are not available. This inference generally gave much larger rates  $\tau^{-1}$  and  $\theta^{-1}$  (Fig. S4A), suggesting confounding effects of the noise on both parameters (reversion to the mean for  $\tau^{-1}$ , and larger variance for  $\theta^{-1}$ ). Indeed, results obtained for a larger value of the frequency threshold ( $f_{th} = 10^{-4}$ , Fig. S4B) gave smaller values, and in better agreement with the age dependence.

To ask whether the clonal dynamics depended on the cell type, we separately analyzed the longitudinally sampled CD4 and CD8 repertoires of P6, the only individual for which such data were available. The clone size distribution of CD4 falls off with a larger exponent than that of CD8, meaning that its largest expanded clones are relatively smaller (Fig. S5 A), as already noted [28]. We then applied the naive inference procedure with

$f_{th} = 10^{-4}$  (since we did not have replicates for the CD4 and CD8 repertoires). The inference (Fig. S5B) reveals that CD4 clones turn over more slowly than CD8 cells (smaller  $\tau^{-1}$ ), but also have much smaller fluctuations in their sizes (smaller  $\theta^{-1}$ ). This result is consistent with a shorter tail of large clones and a larger  $\alpha$  in CD4 than in CD8 (Fig. S5C).

The inference results can be used to predict the persistence of clones, whose turnover has been discussed in the context of aging [28, 29, 35]. For a given individual, we define persistence as the probability that a clone initially observed at frequency  $\geq \hat{f}_{th} = 10^{-5}$  is re-sampled at a later time. This probability strongly depends on the dynamics of turnover of the TCR repertoire, and therefore on the age of the individual, as well as on the time interval between the two samples (Fig. 4D). We can esti-



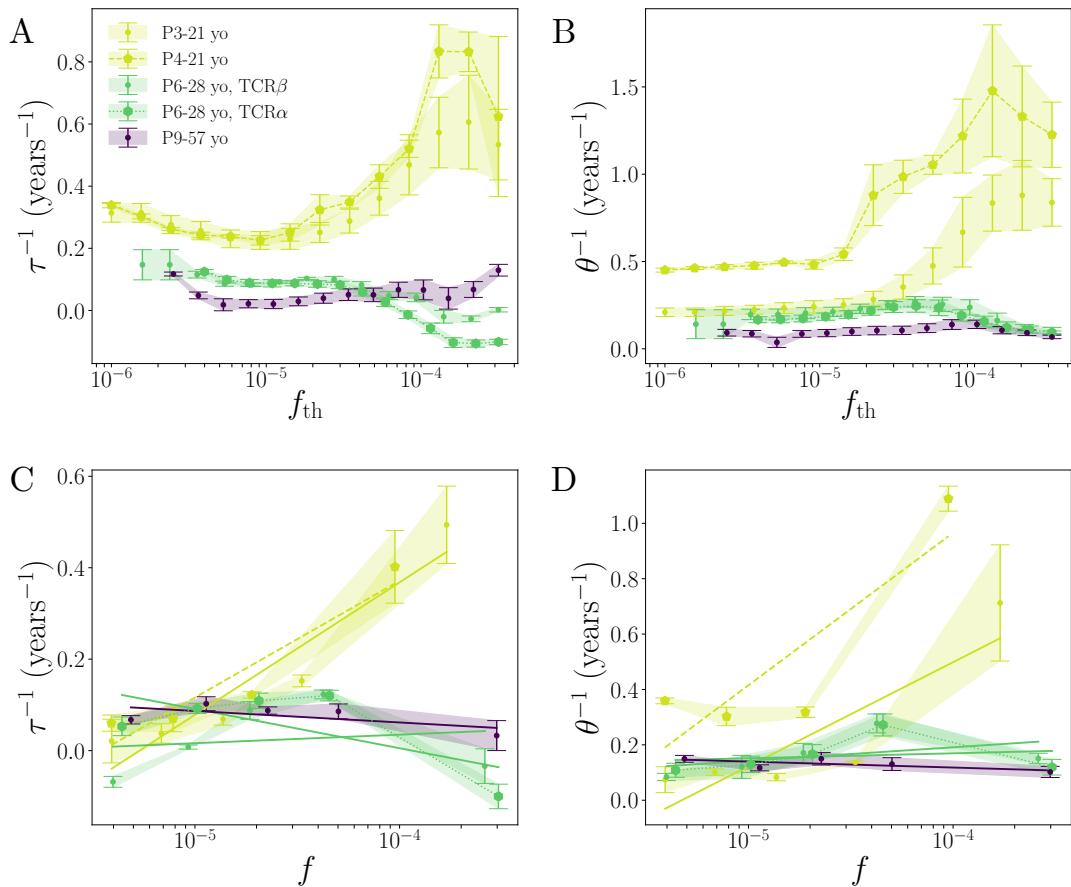
**FIG. 4: Dynamical parameters of healthy TCR repertoires.** **A.** Typical decay rate  $\tau^{-1}$ , and inverse fluctuation amplitude  $\theta^{-1}$  for the 5 donors for whom replicates were available, as obtained using the full inference procedure with  $f_{th} = 10^{-5}$ . All donors but one are consistent with the relation  $2\theta = \tau$ , corresponding to  $\alpha = 1$ . Error bars are standard deviations over all combinations of the replicates at each time point. **B.** Direct test of the prediction  $\alpha = 2\theta/\tau$ . Most values of  $\alpha$  fall close to one, allowing for only a narrow range of tested values. Error bars on  $\alpha$  show standard deviations across time points. **C.** Turnover parameter  $\tau$  as a function of donor age. **D.** Probability for a clone detected at some timepoint with frequency  $\hat{f} \geq 10^{-5}$  to be detected again at a later time point (with the experimental dataset size). Symbols are empirical estimates. The model predictions show excellent agreement. Error bars in **B.-D.** are propagated from **A**.

mate this persistence probability directly from data, and compare it to the predictions of our inferred dynamics, showing excellent agreement. This analysis show that even moderately large clones persist for many years and even decades in older individuals.

Our model assumes that clones have unique trajectories, but that the statistical properties of these trajectories are uniform. However, because of their distinct histories and phenotypical compositions, clones may differ in those dynamical properties. To investigate that possibility, we asked whether the inferred time scales  $\tau$  and  $\theta$  depended on the value of the clone size threshold  $\hat{f}_{th}$ . Low values of  $\hat{f}_{th}$  mean that all clones are taken into account in the inference, while high values mean that we focus on the largest clones only. We found that the values of both  $\tau^{-1}$  and  $\theta^{-1}$  increase with the threshold (Fig. 5A and B) for P3, P4, and P9, suggesting that large clones tend to be more dynamic, with a faster turnover. Therefore, while the overall trends reported in Fig. 4 are

still correct, these results imply that the model should be revisited to allow for frequency-dependent dynamics.

To measure the frequency dependence of the dynamic parameters more finely, we separately inferred  $\tau$  and  $\theta$  for clones sorted into contiguous intervals according to their initial count  $\hat{n}(t_1)$ . Both time scales showed an approximately linear dependency on the logarithm of the initial frequency (Fig. 5 C and D). This confirms the observation that large clones tend to have faster dynamics than small ones, especially in younger individuals. The two time scales  $\tau$  and  $\theta$  vary in concert, so that their ratio remains approximately constant across frequencies (Fig. S6). As we have argued before, this ratio is linked to the power-law exponent of the distribution of clone sizes at steady state. This exponent can be read off as the slope of that distribution on a double logarithmic scale, which it is consistently observed to be constant in the data (Fig. 1D). Finally, we applied the naive inference procedure to learn the frequency-dependent dynamics of



**FIG. 5: Clonal dynamics are frequency dependent.** **A-B.** Results of the full inference as a function of the minimal frequency threshold  $\hat{f}_{\text{th}}$  for  $\tau^{-1}$  and  $\theta^{-1}$ . **C-D.** Dynamical parameters as a function of clone frequency. The inference was performed on separate subsets of clones sorted by their frequency in intervals  $n_{\min} < \hat{n} \leq n_{\max}$ , with  $n_{\min,\max}$  consecutive numbers in  $(2, 5, 10, 20, 100, \infty)$ . Error bars are estimated as in Fig. 4.

clones in all 9 individuals. As expected, this inference yielded more noisy and less stable results than the full inference, especially at low frequencies for which noise is largest (Fig. S7). The dependence of the inferred parameters on clonal frequency vary across individuals, but confirm a picture in which older individuals have more stable large clones.

### III. DISCUSSION

The sizes of T cell clones change constantly throughout the lifetime of an individual, not only due to specific stimulation. We used data sampled on timescales of the order of a year from individuals that did not undergo any strong identified antigenic stimulations to learn the repertoire turnover dynamics. These dynamics include both random unstimulated T cell proliferation and death as well as asymptomatic, or weakly symptomatic antigenic stimulation. We showed that a geometric Brownian motion correctly captures the clone dynamics. This model imposes strict relations that link the exponent of

the TCR clone size distribution at steady state,  $\alpha$ , with the parameters of the dynamics,  $\tau$  and  $\theta$ . We showed that, for all individuals, we were able to predict the measured exponent  $\alpha \approx 1$  from the inferred dynamical parameters, suggesting that geometric Brownian motion is a good description of the process. Although the actual timescales vary between individuals, with much younger individuals having much faster clone turnover dynamics than older individuals, their ratio is fixed. Indeed, as already noted on a larger cohort of individuals in Ref. [35], the exponent of the power-law distribution of clone sizes does not depend on age.

The source of the faster turnover in younger individuals is not explained by our analysis. It can be linked to a larger thymic output rate [1], imposing a faster turnover. It could also be linked to more rapid formation of new immune memories at a young age. We did not attempt to separately learn the dynamics of memory and naive pools, since we did not have sorted longitudinal data for which abundance information could be trusted. While it is sometimes assumed that larger clones have a memory phenotype because they must have expanded, a recent

study in mice has shown that naive clones can be large as well [14]. It will be interesting to perform a separate analysis of careful sorted naive and memory repertoires in the future using the method described here, especially for individuals of different ages.

More generally, we expect clonal dynamics to be linked to the cellular phenotype, as our preliminary analysis showed for CD4 and CD8 cells. Phenotypes can be characterized with increasing resolution using single-cell expression data [37], which also provides paired TCR information [38]. Future work combining longitudinal sampling with single-cell techniques could help explore the relationship between neutral clonal dynamics and cell type. Additionally, we know that TCR with similar sequences form clusters that often respond to similar stimulants [17, 39], and methods are being developed to annotate repertoire with cluster membership [40] or specificity [41–45]. As these annotation become comprehensive, one will be able to study the dynamics of specificity clusters, and to assess the persistence of specific immune memories across different immune challenges.

Our current model is based on two effective parameters that describe the timescale for clone turnover,  $\tau$ , and the timescale of random changes,  $\theta$ . Two major assumptions underlie this model. First, it assumes that antigenic stimulation happens repeatedly on short time scales, so that its cumulated effect on longer time scales look like random fluctuations of the net growth rate. Testing this assumption would require longer time traces of the clonal dynamics, to look for memory effects in the clonal growth rates. Second, it assumes that dynamical properties do not depend on the clone size. As observed in Fig. 5, this assumption is only partially verified, with clear violations for 2 of the youngest donors, in which the larger clones display much faster dynamics than the smaller ones. The longitudinal analysis of larger cohorts with a broad age distribution would be required to investigate this effect in detail.

The turnover time scales we infer range from a few years to 50 years, depending on the age of the individual. It has been shown that even sparsely sampled T-cell repertoires can provide a fingerprint that uniquely identifies individuals [46]. The stability of this immune fingerprint is guaranteed for tens of years, provided that the turnover rate is of the order of years or more, as we showed here.

Direct measurements of T cell lifetimes using heavy water [10] give lifetimes of months for memory cells, to a few years for naive cells. These estimates are consistent with our findings: our time scale  $\tau$  is linked to the inverse of the net growth rate of the *clone*, which results from the balance between cell proliferation and death, while experiments based on heavy water measure the turnover of individual *cells*. For instance, memory cells are short lived, but also divide rapidly to compensate for death, so that the size of memory clones remains stable. One may also want to compare our estimate with the previously reported persistence time of clonotypes believed to be of

fetal origin,  $\approx 37$  years [47]. This persistence time is not directly comparable to  $\tau$ , which is the decay rate of the abundance of each clone, but it is similar to the characteristic decay of the persistence probability (Fig. 4D), which may be slower. Another caveat is that fetal clonotypes are also primarily naive and take up only a few percent of the repertoire, so that they may not be representative of the overall properties of the clonal dynamics.

Our work was possible because we were able to calibrate the noise using replicate samples. However, replicates are not always available. In this case, the dynamics can still be learned for large clones: we showed using simulations that above a certain frequency threshold, the sampling error becomes small and we can use empirical observations to learn TCR repertoire dynamics directly from read counts. This allows us to correctly estimate the dynamics of large clones without a noise model, if the clones sizes are large at both time points. However, since the repertoire is described by a power law distribution, the role of small clones is far from negligible. An alternative to replicates may be to use close-by timepoints (relative to the time scales of the dynamics) as surrogate replicates. While we had such time points separated by one month for P1, P2 and P7, we did not attempt a full inference on these samples: we did not manage learn a reliable noise model for these donors, because we lacked both the raw sequencing reads and details about the processing procedure (PCR amplification, error correction, etc). In particular, unlike uniquely barcoded cDNA sequencing, PCR amplification of gDNA used for these donors inflates rare clonotypes (as suggested by the low-frequency plateau in the clone size distributions, see Fig. 1D), potentially confounding the analysis.

One of the main conclusions of our work is that repertoires are very dynamic systems, with clone frequencies changing by orders of magnitude on timescales of years, even in the absence of strong known stimulation. This observation challenges our ability to identify responding clonotypes to direct immune stimulation, such as vaccination or diseases. This work builds the ground for inference procedures that not only correct for experimental and biological noise but also for the natural repertoire dynamics. The methods we designed are general and can be used on larger cohorts of individuals presenting different health status, age, and immunodeficiencies features. They provide a promising tool to better understand the maintenance and efficiency of T-cells, enabling to quantify immunosenescence [48], which plays an important role in vaccines performance and cancer research.

## IV. METHODS

### A. Longitudinal data

The datasets analyzed in this study are summarized in Table S1, along with accession number and links to databases.

Data was collected from 4 different studies, which uses two different techniques for repertoire sequencing. Data from [28, 30–32] was generated by sequencing TCR mRNA of PBMCs from healthy individuals, while data from [29] was obtained by directly sequencing genomic DNA (gDNA), as described in detail in each original study.

Briefly, mRNA sequencing was done through cDNA synthesis with template switch allowing for the addition of a unique molecular identifier (UMI), followed by 2-step PCR amplification of the TCR loci (alpha and/or beta), multiplexing, sequencing on an Illumina platform, and processing using the MiXCR software package [49], to obtain lists of clonotypes (V and J segments and Complementarity Determining Region 3 nucleotide sequence) corrected for UMI multiplicity and sequencing errors. gDNA sequencing was done by extracting genomic DNA and performing multiplex PCR to amplify the TCR beta gene before sequencing on an Illumina HiSeq system. Raw data processing was performed using closed software. Since the raw data is not available, we used the processed data provided on the ImmuneAccess platform.

## B. Naive inference

The *naive inference* method directly uses the observed TCR clonal frequencies to learn  $\tau$  and  $\theta$  parameters, assuming that they represent exactly the true frequencies:  $f = \hat{f} = \hat{n}/N_r$ . We aim here at maximizing directly the log-likelihood  $\mathcal{L}(\tau, \theta) = \log \mathbb{P}(\{(f_i(t_1), f_i(t_2))\} | \tau, \theta)$ , which can be expressed by integrating Eq.1:

$$\begin{aligned} & \mathbb{P}(\{(\ln f_i(t_1), \ln f_i(t_2))\} | \tau, \theta) \\ &= \prod_i^{N_c} G(\ln f_i(t_2) | \ln f_i(t_1); \tau, \theta) P(\ln f_i(t_1)), \end{aligned} \quad (4)$$

where

$$G(x|y; \tau, \theta) = \sqrt{\frac{\theta}{2\pi\Delta t}} \exp -\frac{\theta(x-y-\Delta t\tau^{-1})^2}{2\Delta t} \quad (5)$$

is the propagator of the Brownian motion,  $\Delta t = t_2 - t_1$  the time interval between the two time points, and where we have assumed that  $N_{\text{cell}}$  is a constant of time. Maximizing the log-likelihood with respect to  $\tau$  and  $\theta$  is equivalent to doing linear regression of  $\ln f(t_2) - \ln f(t_1)$  against  $\Delta t$ .

## C. Full inference

Using same-day replicates at time  $t_j$ , we jointly learn the parameters  $(\alpha, f_{\min}, a, b)$  of the clone-size distribution  $\rho(f) = Cf^{-1-\alpha}$  (for  $f_{\min} \leq f \leq 1$ ), and the noise model  $P(\hat{n}|f) = \text{NegBin}(\hat{n}; N_r f, N_r f + a(N_r f)^b)$  using the NoiSET software [34], where  $\text{NegBin}(n; x, \sigma)$  is a negative binomial of mean  $x$  and variance  $\sigma$ . The learned parameters are reported in Fig. S1.

We then learn parameters of the dynamics by maximizing the likelihood of samples taken at two different time points, using the noise model to account for the discrepancy between true frequencies and sequence counts: For one clone, the full model likelihood reads

$$\begin{aligned} & \mathbb{P}(\hat{n}_i(t_1) = \hat{n}_1, \hat{n}_i(t_2) = \hat{n}_2) | \tau, \theta = \\ & \int_{[f_{\min}, 1]^2} df_1 \rho(f_1) \frac{df_2}{f_2} G(\ln f_2 | \ln f_1; \tau, \theta) P(\hat{n}_1 | f_1) P(\hat{n}_2 | f_2). \end{aligned} \quad (6)$$

where the noise models are specific to each time point.

The maximum likelihood estimator is given by :

$$(\tau^*, \theta^*) = \underset{(\tau, \theta)}{\operatorname{argmax}} \prod_{i=1}^{N_c} \frac{\mathbb{P}(\hat{n}_i(t_1), \hat{n}_i(t_2) | \tau, \theta)}{\mathbb{P}(\hat{n}_i(t_1) \geq N_r f_{\text{th}}, \hat{n}_i(t_2) > 0 | \tau, \theta)}, \quad (7)$$

where the denominator accounts for the condition that the clone be included in the analysis:  $\hat{f}_i(t) \geq f_{\text{th}}$  and  $\hat{n}_2 > 0$ . The persistence probability of Fig. 4D is linked to that normalization and is computed as:

$$P_{\text{pers}}(\tau, \theta) = \frac{\mathbb{P}(\hat{n}_i(t_1) \geq N_r f_{\text{th}}, \hat{n}_i(t_2) > 0 | \tau, \theta)}{\mathbb{P}(\hat{n}_i(t_1) \geq N_r f_{\text{th}} | \tau, \theta)}. \quad (8)$$

Once the model is learned, the posterior distribution of fold changes  $s_i \equiv \ln f_i(t_2) - \ln f_i(t_1)$  of each clone  $i$  is computed through

$$\begin{aligned} & \mathbb{P}(s_i = s | \hat{n}_1, \hat{n}_2, \tau^*, \theta^*) = \\ & \frac{\int_{f_{\min}}^1 df_1 \rho(f_1) G(\ln f_1 + s | \ln f_1; \tau^*, \theta^*) P(\hat{n}_1 | f_1) P(\hat{n}_2 | f_1 e^s)}{\mathbb{P}(\hat{n}_1, \hat{n}_2 | \tau^*, \theta^*)}. \end{aligned} \quad (9)$$

The overall posterior distribution over all clones (solid lines in Fig. S3) is then given by  $\mathbb{P}_{\text{post}}(s) = (1/N_c) \sum_{i=1}^{N_c} \mathbb{P}(s_i = s | \hat{n}_1, \hat{n}_2, \tau^*, \theta^*)$ . The prior distribution (dashed line), by contrast, is directly given by  $G(\ln f_1 + s, \ln f_1 | \tau^*, \theta^*)$ , which is independent of  $f_1$ .

When doing inference in each frequency bin, the product in (7) runs over clones that fall in the bin, and the normalization in the denominator is replaced by the probability to observe  $\hat{n}_i(t_1)$  in the bin of interest, and  $\hat{n}_i(t_2) > 0$ . The maximization is performed using the `minimize` function from the Scipy package, with the Sequential Least Squares Programming (SLSQP) method [50] with parameters `tol=1e-8` and `maxiter=300` and initial condition  $\tau = 2, \theta = .5$  and constraint  $\theta^{-1} > 10^{-3}$ .

## D. Synthetic data

Synthetic data was generated by simulating Eq. 1 with a source term producing new clones with rate  $S$  at initial size  $n = n_0 = 40$ , and an absorbing boundary condition at  $n = 1$ . We work with the  $x = \ln n$  variable for convenience. The simulation is initialized at

steady state, which can be computed analytically [5, 9]. The analytical solution gives us the expected number of cells and clones as a function of the model parameters:  $N_{\text{cell}} = S(n_0 - 1)/(\tau^{-1} - \theta^{-1}/2)$ , and  $N_c = S\tau \ln n_0$ . Fixing the number of cells to  $N_{\text{cell}} = 10^{10}$ , we then compute the number of clones necessary to achieve that size,  $N_c = N_{\text{cell}}(1 - \tau/2\theta) \ln n_0/(n_0 - 1)$ . We then draw the size  $n_i(t_1) = e^{x_i(t_1)}$  of each clone  $i = 1, \dots, N_c$  from the continuous steady state distribution [5]:

$$\rho_x(x) = \begin{cases} S\tau(1 - e^{-\alpha x}) & \text{if } x \leq x_0 \equiv \ln n_0 \\ S\tau e^{-\alpha x}(e^{\alpha x_0} - 1) & \text{if } x > x_0, \end{cases} \quad (10)$$

with  $\alpha = 2\theta/\tau$ .

Then the evolution of each clone from time  $t_1$  to  $t_2 = t_1 + \Delta t$  is determined by the modified propagator with absorbing boundary condition at  $x = 0$ :

$$G_{\text{abs}}(x|y) = G(x|y) - e^{-\alpha y}G(x| - y) \quad (11)$$

where  $G(x|y)$  is defined in (5). In practice, we kill clone  $i$  with probability  $1 - P_{\text{surv}}(x_i(t_1)) \equiv 1 - \int_0^\infty dx G_{\text{abs}}(x|x_i(t_1))$ , which can be expressed in terms of error functions. Otherwise, its new log-size  $x_i(t_2)$  is drawn from the distribution  $G_{\text{abs}}(x|x_i(t_1))/P_{\text{surv}}(x_i(t_1))$ .

In addition, new clones are introduced during  $\Delta t$ . We draw their number from a Poisson distribution of mean

$S\Delta t$ , and their introduction times  $t$  from a uniform distribution in the interval  $[t_1, t_2]$ . Then their dynamics is drawn in the same way as for the initial clones, but with  $\Delta t = t_2 - t$  instead of  $t_2 - t_1$ .

Once the abundances ( $n_i(t_1) = e^{x_i(t_1)}$ ,  $n_i(t_2) = e^{x_i(t_2)}$ ) have been determined, the number of reads  $\hat{n}_i(t_1)$ ,  $\hat{n}_i(t_2)$  from each time point is drawn from a negative binomial distribution of mean  $\langle \hat{n}_i(t_1) \rangle = N_r n_i(t_1)/N_{\text{cell}}$  and variance  $\langle \hat{n}_i(t_1) \rangle + a \langle \hat{n}_i(t_1) \rangle^b$ , and likewise for  $\hat{n}_i(t_2)$ , with  $N_r = 10^6$ ,  $a = 0.7$  and  $b = 1.1$ .

## E. Code availability

All scripts to produce the figures can be found at [https://github.com/statbiophys/Inferring\\_TCR\\_reertoire\\_dynamics/](https://github.com/statbiophys/Inferring_TCR_reertoire_dynamics/).

## Acknowledgements

This work was partially supported by the European Research Council Consolidator Grant n. 724208 and ANR-19-CE45-0018 “RESP-REP” from the Agence Nationale de la Recherche.

- 
- [1] Yates AJ (2014) Theories and Quantification of Thymic Selection. *Frontiers in Immunology* 5.
  - [2] Bains I, Thiébaut R, Yates AJ, Callard R (2009) Quantifying Thymic Export: Combining Models of Naive T Cell Proliferation and TCR Excision Circle Dynamics Gives an Explicit Measure of Thymic Output. *The Journal of Immunology* 183:4329–4336.
  - [3] Jameson SC (2002) Maintaining the norm: T-cell homeostasis. *Nature Reviews Immunology* 2:547–556.
  - [4] Dowling MR, Hodgkin PD (2009) Modelling naive T-cell homeostasis: Consequences of heritable cellular lifespan during ageing. *Immunology & Cell Biology* 87:445–456.
  - [5] Desponts J, Mora T, Walczak AM (2016) Fluctuating fitness shapes the clone-size distribution of immune repertoires. *Proceedings of the National Academy of Sciences* 113:274–279.
  - [6] Mora T, Walczak AM (2019) in *Systems Immunology* (CRC Press).
  - [7] Qi Q, et al. (2014) Diversity and clonal selection in the human T-cell repertoire. *Proceedings of the National Academy of Sciences* 111:13139–13144.
  - [8] Lythe G, Callard RE, Hoare RL, Molina-Paris C (2016) How many TCR clonotypes does a body maintain? *Journal of Theoretical Biology* 389:214–224.
  - [9] Mora T, Walczak AM (2019) How many different clonotypes do immune repertoires contain? *Current Opinion in Systems Biology* 18:104–110.
  - [10] De Boer RJ, Homann D, Perelson AS (2003) Different Dynamics of CD4<sup>+</sup> and CD8<sup>+</sup> T Cell Responses During and After Acute Lymphocytic Choriomeningitis Virus Infection. *The Journal of Immunology* 171:3928–3935.
  - [11] Kedzierska K, Valkenburg SA, Doherty PC, Davenport MP, Venturi V (2012) Use it or lose it: Establishment and persistence of T cell memory. *Frontiers in Immunology* 3.
  - [12] Mayer A, Zhang Y, Perelson AS, Wingreen NS (2019) Regulation of T cell expansion by antigen presentation dynamics. *Proceedings of the National Academy of Sciences* 116:5914–5919.
  - [13] Bains I, Antia R, Callard R, Yates AJ (2009) Quantifying the development of the peripheral naive CD4+ T-cell pool in humans. *Blood* 113:5480–5487.
  - [14] de Greef PC, et al. (2020) The naive T-cell receptor repertoire has an extremely broad distribution of clone sizes. *eLife* 9:e49900.
  - [15] Dessalles R, et al. (2022) How Naive T-Cell Clone Counts Are Shaped By Heterogeneous Thymic Output and Homeostatic Proliferation. *Frontiers in Immunology* 12.
  - [16] De Boer RJ, Perelson AS (1994) T Cell Repertoires and Competitive Exclusion. *Journal of Theoretical Biology* 169:375–390.
  - [17] Dash P, et al. (2017) Quantifiable predictive features define epitope-specific T cell receptor repertoires. *Nature* 547:89–93.
  - [18] Mayer A, Balasubramanian V, Walczak AM, Mora T (2019) How a well-adapting immune system remembers. *Proceedings of the National Academy of Sciences* 116:8815–8823.
  - [19] Johnson PLF, Yates AJ, Goronzy JJ, Antia R (2012) Pe-

- ripheral selection rather than thymic involution explains sudden contraction in naive CD4 T-cell diversity with age. *Proceedings of the National Academy of Sciences* 109:21432–21437.
- [20] De Boer RJ, Perelson AS (2013) Quantifying T lymphocyte turnover. *Journal of Theoretical Biology* 327:45–87.
- [21] Lindau P, Robins HS (2017) Advances and applications of immune receptor sequencing in systems immunology. *Current Opinion in Systems Biology* 1:62–68.
- [22] Davis MM, Boyd SD (2019) Recent progress in the analysis of  $A\beta$  T cell and B cell receptor repertoires. *Current Opinion in Immunology* 59:109–114.
- [23] Minervina A, Pogorelyy M, Mamedov I (2019) T-cell receptor and B-cell receptor repertoire profiling in adaptive immunity. *Transplant International* 32:1111–1123.
- [24] Burgos JD, Moreno-Tovar P (1996) Zipf-scaling behavior in the immune system. *Biosystems* 39:227–232.
- [25] Koch H, Starenki D, Cooper SJ, Myers RM, Li Q (2018) powerTCR: A model-based approach to comparative analysis of the clone size distribution of the T cell receptor repertoire. *PLOS Computational Biology* 14:e1006571.
- [26] Naumov YN, Naumova EN, Hogan KT, Selin LK, Gorski J (2003) A Fractal Clonotype Distribution in the CD8+ Memory T Cell Repertoire Could Optimize Potential for Immune Responses. *The Journal of Immunology* 170:3994–4001.
- [27] Zarnitsyna VI, Evavold BD, Schoettle LN, Blattman JN, Antia R (2013) Estimating the Diversity, Completeness, and Cross-Reactivity of the T Cell Repertoire. *Frontiers in Immunology* 4.
- [28] Britanova OV, et al. (2016) Dynamics of Individual T Cell Repertoires: From Cord Blood to Centenarians. *The Journal of Immunology* 196:5005–5013.
- [29] Chu ND, et al. (2019) Longitudinal immunosequencing in healthy people reveals persistent T cell receptors rich in highly public receptors. *BMC Immunology* 20:19.
- [30] Sycheva AL, et al. (2018) Quantitative profiling reveals minor changes of T cell receptor repertoire in response to subunit inactivated influenza vaccine. *Vaccine* 36:1599–1605.
- [31] Pogorelyy MV, et al. (2018) Precise tracking of vaccine-responding T cell clones reveals convergent and personalized response in identical twins. *Proceedings of the National Academy of Sciences* 115:12704–12709.
- [32] Minervina AA, et al. (2020) Primary and secondary anti-viral response captured by the dynamics and phenotype of individual T cell clones. *eLife* 9:e53704.
- [33] Puelma Touzel M, Walczak AM, Mora T (2020) Inferring the immune response from repertoire sequencing. *PLOS Computational Biology* 16:e1007873.
- [34] Bensouda Koraichi M, Touzel MP, Mora T, Walczak AM (2021) NoisET: Noise learning and Expansion detection of T-cell receptors with Python. *arXiv:2102.03568 [q-bio]*.
- [35] Gaimann MU, Nguyen M, Desponts J, Mayer A (2020) Early life imprints the hierarchy of T cell clone sizes. *eLife* 9:e61639.
- [36] Altan-Bonnet G, Mora T, Walczak AM (2020) Quantitative Immunology for Physicists. *Physics Reports* 849:1–83.
- [37] Pai JA, Satpathy AT (2021) High-throughput and single-cell T cell receptor sequencing technologies. *Nature Methods* 18:881–892.
- [38] Valkiers S, et al. (2022) Recent advances in T-cell receptor repertoire analysis: Bridging the gap with multimodal single-cell RNA sequencing. *ImmunoInformatics* 5:100009.
- [39] Glanville J, et al. (2017) Identifying specificity groups in the T cell receptor repertoire. *Nature* 547:94–98.
- [40] Mayer-Blackwell K, et al. (2021) TCR meta-clonotypes for biomarker discovery with tcrdist3 enabled identification of public, HLA-restricted clusters of SARS-CoV-2 TCRs. *eLife* 10:e68605.
- [41] Gielis S, et al. (2019) Detection of Enriched T Cell Epitope Specificity in Full T Cell Receptor Sequence Repertoires. *Frontiers in Immunology* 10.
- [42] Montemurro A, et al. (2021) NetTCR-2.0 enables accurate prediction of TCR-peptide binding by using paired TCR $\alpha$  and  $\beta$  sequence data. *Communications Biology* 4:1–13.
- [43] Sidhom JW, Larman HB, Pardoll DM, Baras AS (2021) DeepTCR is a deep learning framework for revealing sequence concepts within T-cell repertoires. *Nature Communications* 12:1605.
- [44] Springer I, Tickotsky N, Louzoun Y (2021) Contribution of T Cell Receptor Alpha and Beta CDR3, MHC Typing, V and J Genes to Peptide Binding Prediction. *Frontiers in Immunology* 12.
- [45] Zhang W, et al. (2021) A framework for highly multiplexed dextramer mapping and prediction of T cell receptor sequences to antigen specificity. *Science Advances* 7:eabf5835.
- [46] Dupic T, et al. (2021) Immune fingerprinting through repertoire similarity. *PLOS Genetics* 17:e1009301.
- [47] Pogorelyy MV, et al. (2017) Persisting fetal clonotypes influence the structure and overlap of adult human T cell receptor repertoires. *PLOS Computational Biology* 13:e1005572.
- [48] Zhang H, Weyand CM, Goronzy JJ (2021) Hallmarks of the aging T-cell system. *The FEBS Journal* 288:7123–7142.
- [49] Bolotin DA, et al. (2015) MiXCR: Software for comprehensive adaptive immunity profiling. *Nature Methods* 12:380–381.
- [50] Virtanen P, et al. (2020) SciPy 1.0: Fundamental algorithms for scientific computing in Python. *Nature Methods* 17:261–272.

### Supplementary information

ID	age	sex	# clones	# reads	Tech	DOI ref	data access
P1	18-19 yo	female	4.08 - 5.11·10 <sup>5</sup>	10.8 - 25.1·10 <sup>6</sup>	gDNA	[29]	<a href="https://doi.org/10.21417/B7J01X">https://doi.org/10.21417/B7J01X</a>
P2	18-19 yo	female	1.55 - 2.93·10 <sup>5</sup>	10.5 - 20.7·10 <sup>6</sup>	gDNA	[29]	<a href="https://doi.org/10.21417/B7J01X">https://doi.org/10.21417/B7J01X</a>
P3	21-23 yo	male	2.04 - 6.44·10 <sup>5</sup>	0.23 - 1.03·10 <sup>6</sup>	RNA	[31]	PRJNA493983
P4	21-23 yo	male	1.9 - 10.06·10 <sup>5</sup>	0.28 - 1.79·10 <sup>6</sup>	RNA	[31]	PRJNA493983
P5	27-30 yo	male	7.6 - 18.12·10 <sup>5</sup>	1.53 - 6.94·10 <sup>6</sup>	RNA	[28]	PRJNA316572
P6	28-29 yo	male	2.62 - 6.38·10 <sup>5</sup>	0.56 - 1.69·10 <sup>6</sup>	RNA	[32]	PRJNA577794
P7	45-46 yo	female	1.93 - 2.42·10 <sup>5</sup>	8.86 - 22.8·10 <sup>6</sup>	gDNA	[29]	<a href="https://doi.org/10.21417/B7J01X">https://doi.org/10.21417/B7J01X</a>
P8	47-50 yo	male	1.38 - 9.54·10 <sup>5</sup>	1.53 - 5.93·10 <sup>6</sup>	RNA	[28]	PRJNA316572
P9	57-58 yo	male	3.25 - 7.29·10 <sup>5</sup>	0.62 - 1.64·10 <sup>6</sup>	RNA	[30]	SRP111073

TABLE 1: Summary of the repertoire samples and individuals used in this study.

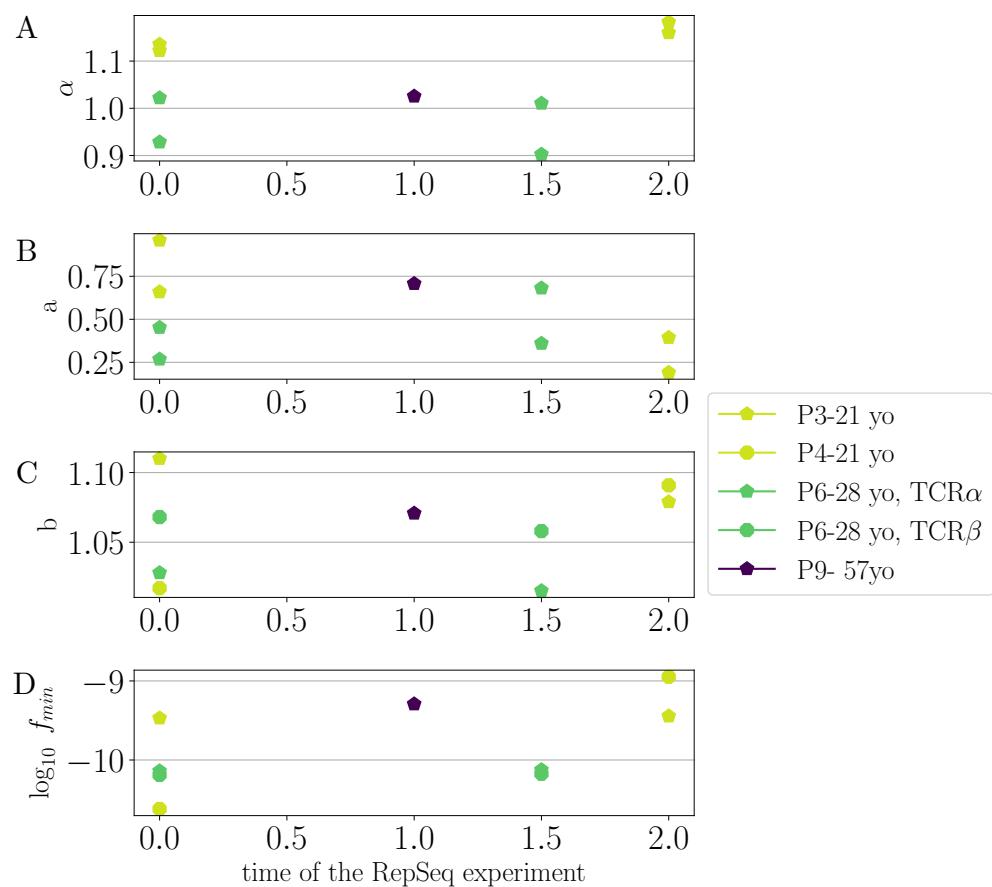


FIG. S1: Parameters of the null model, which characterize both the noise model (*a* and *b*) and the prior power law distribution of frequencies ( $\alpha, f_{\min}$ ). The x-axis represents time in years since the first sample was taken for each individual. Only samples for which duplicates are available are shown.

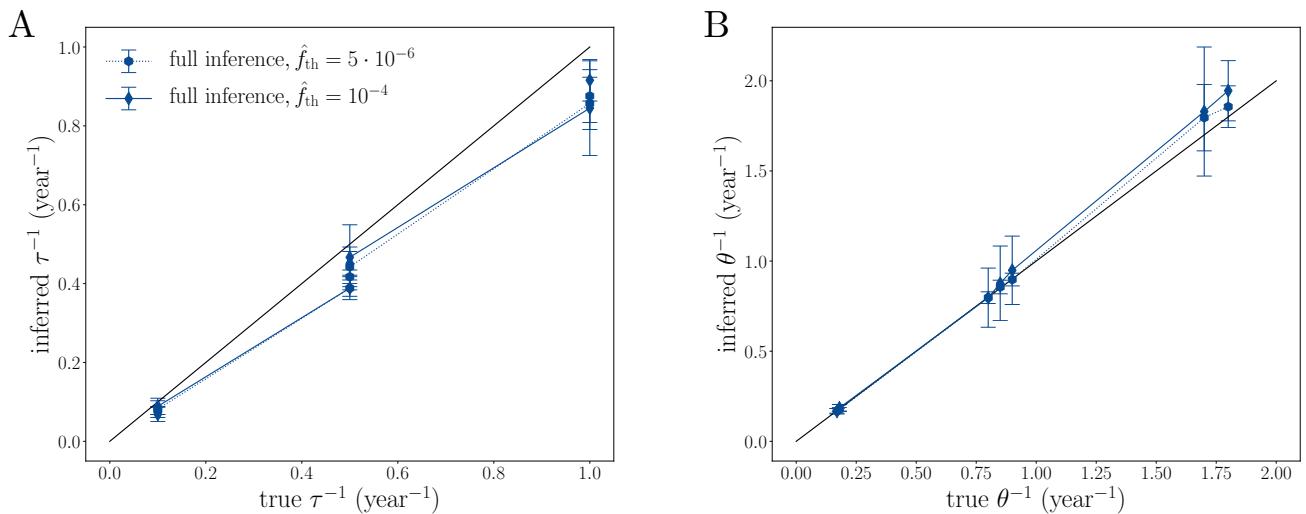


FIG. S2: Validation of the inference method for time point with very different number of reads. Parameters: all 9 combinations of  $\tau^{-1} = (0.1, 0.5, 1)$  year $^{-1}$  and  $\alpha = (1.11, 1.17, 1.25)$ . The number of reads are  $N_r(t_1) = 10^6$  and  $N_r(t_2) = 10^5$ .

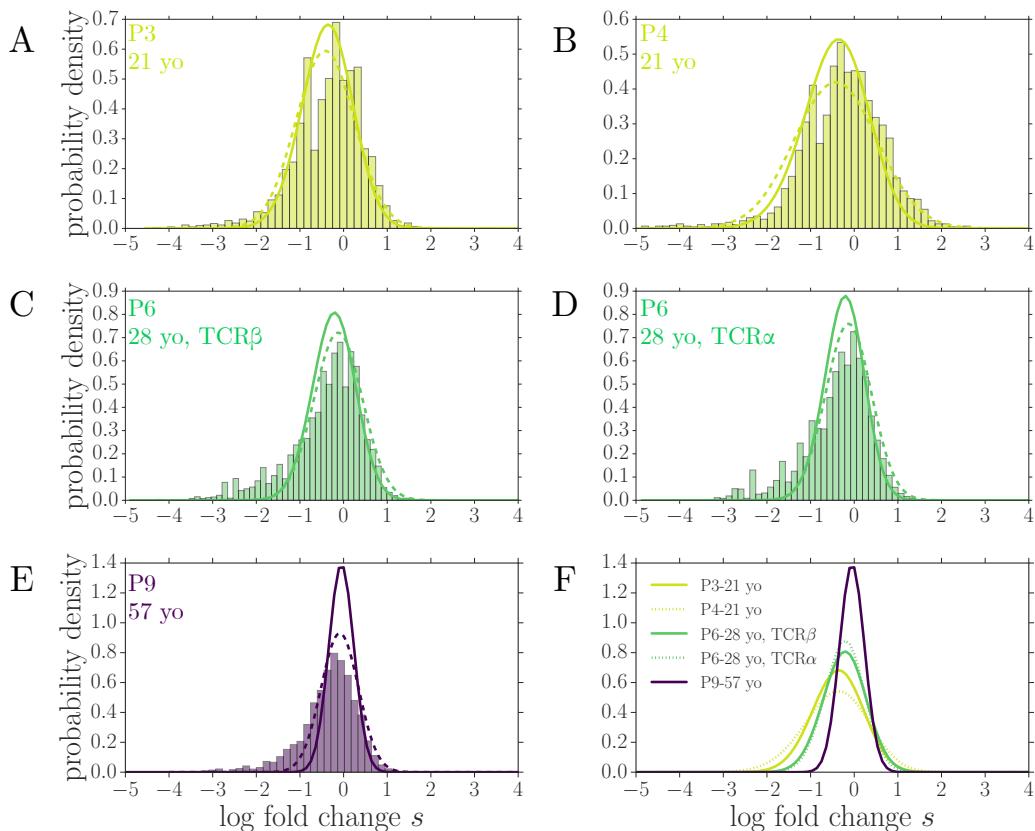


FIG. S3: Distribution of log-fold changes  $s$ . A-E. For each individual and chain we compare the naive distribution  $s \approx \ln(\hat{f}_2/\hat{f}_1)$  (histogram bars), the prior distribution  $G(\ln f_1 + s, \ln f_1 | \tau^*, \theta^*)$  with the inferred parameters, and the posterior distribution  $(1/N_c) \sum_i \mathbb{P}(s | \hat{n}_i(t_1), \hat{n}_i(t_2); \tau^*, \theta^*)$ . While the naive distribution has an excess of low values due to small number errors, the prior and posterior distribution agree well. F. Comparison of posteriors across individuals, showing how both the average decay and its spread decrease with age.

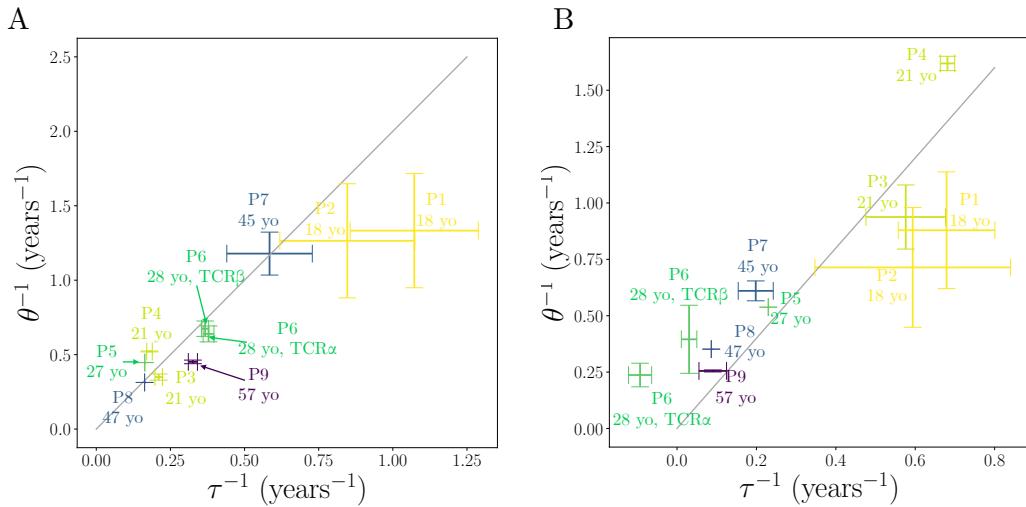


FIG. S4: Naive inference of the dynamical parameters on all individuals, with (A)  $f_{th} = 10^{-5}$  and (B)  $f_{th} = 10^{-4}$ .

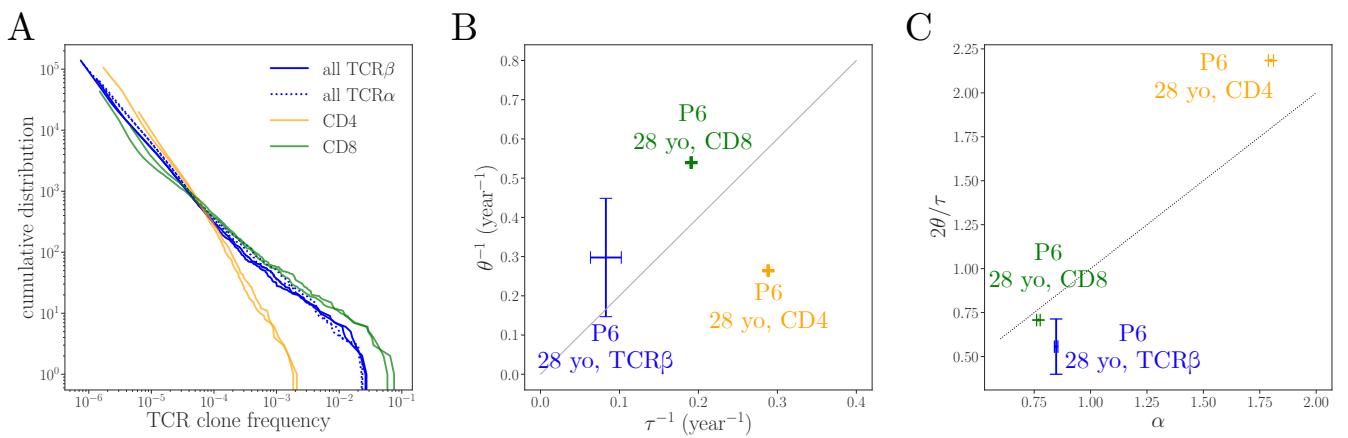


FIG. S5: Comparison for the CD4 and CD8 repertoire dynamics in P6. **A.** Clone size distribution in the bulk (alpha and beta chains) and in the CD4 and CD8 beta chain repertoires. The CD4 repertoire has a shorter tail, corresponding to a larger exponent  $\alpha$ . **B.** Prior and posterior distributions of the log-fold change, as in Fig. S3. **C.** Inferred parameters for each repertoire. **D.** Model prediction ( $2\theta/\tau$ ) vs measured power-law exponent  $\alpha$ . The smaller amplitude of frequency fluctuations  $\theta^{-1}$  in the CD4 repertoire is consistent with its shorter tail of large clones.

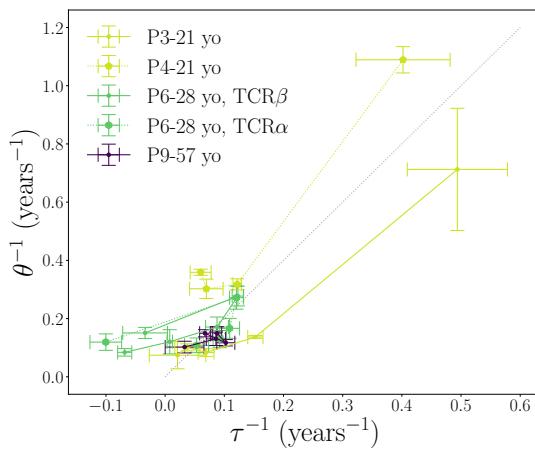


FIG. S6: Inferred values of  $\tau$  and  $\theta$  for different individuals and frequencies. Frequency intervals and datapoints are the same as in Fig. 5 C-D.

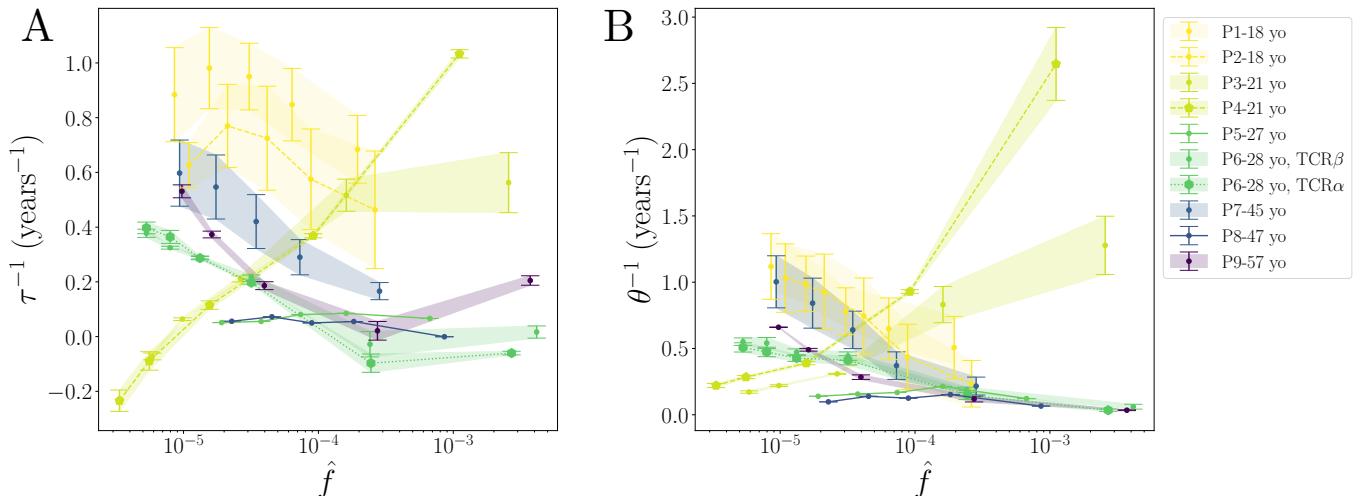


FIG. S7: The inference was performed on separate subsets of clones sorted by their frequency in intervals  $n_{\min} < \hat{n} \leq n_{\max}$ , with  $n_{\min, \max}$  consecutive numbers in  $(3, 5, 7, 15, 10, 1000, \infty)$  for P3, P4, P6, P9, and  $(100, 200, 400, 800, 2000, \infty)$  for P1, P2, P5, P7 and P8.